

4. an audio pickup device for generating audio signals representative of sound from an
5 audio source; and

6 a multimodal integration architecture system for processing said image signals and said
7 audio signals to determine a direction of the audio source relative to a reference point.

1 2. (original) The video conferencing system of claim 1 wherein said multimodal integration
2 architecture system further comprises:

3 an audio source localization system;

4 a computer vision person detection system; and

5 a multimodal speaker detection system.

C1
1 3. (original) The video conferencing system of claim 2, further comprising an integrated
2 housing for an integrated video conferencing system incorporating the image pickup device, the
3 audio pickup device, and the multimodal integration architecture system.

2 4. (original) The video conferencing system of claim 3, wherein the integrated housing is
sized for being portable.

1 5. (original) The video conferencing system of claim 2, further comprising an electronic pan
2 tilt zoom system for electronically manipulating the image signals to effectively provide at least
3 one of variable pan, tilt, and zoom functions.

1 6. (original) The video conferencing system of claim 5, wherein the image pickup device is a
2 stationary camera.

1 7. (original) The video conferencing system of claim 5, wherein the multimodal integrated
2 architecture system provides control signals to the electronic pan tilt zoom system.

09/822,121

3

1 8. (original) The video conferencing system of claim 7, wherein the audio source moves
2 relative to the reference point, the audio source localization system detects the movement of the
3 audio source, and, in response to the movement, the audio source localization system causes a
4 change in the field of view of the image pickup device.

1 9. (original) The video conferencing system of claim 5, wherein the audio pickup device is
2 comprised of an array of two microphones.

1 10. (currently amended) A method comprising the steps of:

2 generating, at a stationary an image pickup device, remaining motionless during
3 operation, image signals representative of an image;

4 generating, at an audio pickup device, audio signals representative of sound from an
5 audio source;

6 processing the image signals and the audio signals to determine a direction of the audio
7 source relative to a reference point;

8 manipulating the image signals to produce refined image signals; and

9 outputting said refined image signals.

1 11. (original) The method of claim 10 further comprising the steps of:

2 applying said audio signals to an audio source localization system;

3 applying said image signals to a computer vision person detection system;

4 processing said audio signals and said image signals with a multimodal speaker detection
5 system;

6 generating control signals based on the determined direction of the audio source;

7 applying the control signals to an electronic pan tilt zoom system to mimic the effect of at
8 least one function of a movable camera, said function selected from the group consisting
9 panning, tilting, and zooming said movable camera; and

09/822,121

4

10 providing an output from said electronic pan tilt zoom system.

1 12. (original) The method of claim 10, further comprising electronically varying a field of
2 view of the image pickup device in response to the control signals.

1 13. (original) The method of claim 10, wherein processing the audio signals includes
2 determining an audio based direction of the audio source based on the audio signals.

1 14. (original) The method of claim 12, wherein the audio source moves relative to a
2 reference point, and wherein processing the audio signals further includes:

3 detecting the movement of the audio source; and

4 causing electronically, in response to the movement, an increase in the field of view of
5 the image pickup device.

1 15. (original) The method of claim 12, further comprising the step of supplying control
2 signals, based on the audio based direction, for electronically panning, tilting, or zooming said
3 image pickup device.

1 16. (currently amended) A video conferencing system comprising:
2 two microphones for generating audio signals representative of sound from a speaker;
3 a stationary video camera, remaining motionless during operation, for generating video
4 signals representative of a video image;
5 an electronic pan tilt zoom system for manipulating video images to produce the visual
6 effects of panning, tilting, and/or zooming;
7 a processor for processing the video signals and the audio signals to determine a direction
8 of a speaker relative to a reference point and supplying control signals to the electronic pan tilt
9 zoom system for producing images that include the speaker in the field of view of the camera, the

09/822,121

5

10

control signals being generated based on the determined direction of the speaker; and

11

a transmitter for transmitting audio and video signals for video conferencing.